

January 2013

Building a Bilingual Corpus based on Hybrid Approach for Malayalam-English Machine Translation

Rajesh. K. S

CMRCET, Hyderabad, India., ktmksrajesh@gmail.com

Veena A Kumar

MGIT, Hyderabad, India, veenaakumar@gmail.com

CH. Dayakar Reddy

CMRCET, Hyderabad, India, daya_chintla@yahoo.co.in

Follow this and additional works at: <https://www.interscience.in/ijcsi>



Part of the [Computer Engineering Commons](#), [Information Security Commons](#), and the [Systems and Communications Commons](#)

Recommended Citation

S, Rajesh. K.; Kumar, Veena A; and Reddy, CH. Dayakar (2013) "Building a Bilingual Corpus based on Hybrid Approach for Malayalam-English Machine Translation," *International Journal of Computer Science and Informatics*: Vol. 2 : Iss. 3 , Article 13.

Available at: <https://www.interscience.in/ijcsi/vol2/iss3/13>

This Article is brought to you for free and open access by Interscience Research Network. It has been accepted for inclusion in International Journal of Computer Science and Informatics by an authorized editor of Interscience Research Network. For more information, please contact sritampatnaik@gmail.com.

Building a Bilingual Corpus based on Hybrid Approach for Malayalam-English Machine Translation

Rajesh. K. S¹, Veena A Kumar² & CH. Dayakar Reddy³

^{1,3} CMRCET, Hyderabad, India, ²MGIT, Hyderabad, India
E-mail : ktmksrajesh@gmail.com¹, veenaakumar@gmail.com²,
daya_chintla@yahoo.co.in³

Abstract - Word alignment in bilingual corpora has been a very active research topic in the Machine Translation research groups. In this research paper, we describe an alignment system that aligns English-Malayalam texts at word level in parallel sentences. The alignment of translated segments with source segments is very essential for building parallel corpora. Since word alignment research on Malayalam and English languages is still in its immaturity, it is not a trivial task for Malayalam-English text. A parallel corpus is a collection of texts in two languages, one of which is the translation equivalent of the other. Thus, the main purpose of this system is to construct word-aligned parallel corpus to be used in Malayalam-English machine translation. The proposed approach is a hybrid approach, a combination of corpus based and dictionary lookup approaches. The corpus based approach is based on the first three IBM models and Expectation Maximization (EM) algorithm. For the dictionary lookup approach, the proposed system uses the bilingual Malayalam-English Dictionary.

Key words - EM Algorithm, IBM Models, Machine Translation, Word-aligned Parallel Corpus, Natural Language Processing.

I. INTRODUCTION

Processing the Malayalam texts is difficult in its computation because sentences in Malayalam texts are represented as strings of Malayalam characters without spaces to indicate word boundaries. This cause various problems for Machine Translation, Information Retrieval, Text Summarization and many other Natural Language Processing activities. Bilingual word alignment is the first step of most current approaches to Statistical Machine Translation or SMT. One simple and very old but still quite useful approach for language modeling is n-gram modeling. Separate language models are built for the source language (SL) and the target language (TL). For the first stage, monolingual corpora of the SL and the TL are required. The second stage is called translation modeling which includes the step of finding the word alignments induced over a sentence aligned bilingual (parallel) corpus. This paper particularly deals with the step of word alignment. Corpora and other lexical resources are not yet widely available in Malayalam and hence research in language technologies has not progressed much. In this paper we describe our efforts in building an English-Malayalam aligned parallel corpus. Although parallel corpora are very useful resources for many natural languages processing applications such as building machine translation systems, multilingual dictionaries and word sense disambiguation, they are not yet available for

many languages of the world, particularly Indian languages like Malayalam. Building a parallel corpus manually is a very tedious and time-consuming task. A good way to develop such a corpus is to start from available resources containing the translations from the source language to the target language. A parallel corpus becomes very useful when the texts in the two languages are aligned. This system uses the IBM models to align the texts at word level.

Many words in natural languages like Malayalam have multiple meanings. It is important to identify the correct sense of a word before we start translation, query-based information retrieval, information extraction, question answering, etc. Recently, parallel corpora are being employed for detecting the correct sense of a word. Ng [2] proposed that if two languages are not closely related, different senses in the source language are likely to be translated differently in the target language. Parallel corpus based techniques for word sense disambiguation work better when the two languages are dissimilar.

The paper is structured as follows. Section 2 describes some related work. Section 3 presents the Alignment Model. Section 4, discusses the Proposed Alignment Model. Section 5 gives an Overview of System. In section 6, the expected experimental results are presented. Lastly, section 7 presents conclusion and future work.

II. RELATED WORK

An enormous amount of research has been conducted in the alignment of parallel texts with various methodologies. G. Chinnappa and Anil Kumar Singh [6] proposed a java implementation of an extended word alignment algorithm based on the IBM models. By introducing a similarity measure (Dice coefficient), using a list of cognates and morph analyzer, they have been able to improve the performance. Li and Chengqing Zong [11] addressed the word alignment between sentences with different valid word orders, which changes the order of the word sequences (called word reordering) of the output hypotheses to make the word order more exactly match the alignment reference.

Pascale Fung and Kenneth Ward Church, in K-vec algorithm [13], makes use of the word position and frequency feature to find word correspondences using Euclidean distance. Ittycheriah and Roukos [8] proposed a maximum entropy word aligner for Arabic-English machine translation. Martin et al. [9] have discussed word alignment for languages with scarce resources. Bing Xiang, Yonggang Deng and Bowen Zhou [1] proposed Diversify and Combine: Improving Word Alignment for Machine Translation on Low-Resource Languages. This approach on an English-to-Pashto translation task by combining the alignments obtained from syntactic reordering, stemming, and partial words. Jamie Brunning, Adria de Gispert and William Byrne proposed Context-Dependent Alignment Models for Statistical Machine Translation [10]. These models lead to an improvement in alignment quality, and an increase in translation quality when the alignments are used in Arabic-English and Chinese-English translation.

Most current SMT systems [14] presently available use a generative model for word alignment such as the one implemented in the freely available tool GIZA++ [16]. GIZA++ is an implementation of the IBM alignment models [15]. These models treat word alignment as a hidden process, and maximize the probability of the observed (e, f) sentence pairs using the Expectation Maximization (EM) algorithm, where e and f are the source and the target sentences. In [4] all the conducted experiments prove that the augmented approach, on multiple corpora, performs better when compared to the use of GIZA++ and NATools individually for the task of English-Hindi word alignment. D.Wu, (1994) [3] has developed Chinese and English parallel corpora in the Department of Computer Science and University of Science and Technology, Hong Kong. Here two methods are applied which are important. Firstly, the Gale's methods is used to Chinese and English which shows that length-based methods give satisfactory result even between unrelated languages which is a surprising result. Next, it shows the

effect of adding lexical cues to length-based methods. According to these results, using lexical information increases accuracy of alignment from 86% to 92%.

A hybrid approach to align sentences and words in English-Hindi parallel corpora [12] presented an alignment system that aligns English-Hindi texts at the sentence and word levels in parallel corpora. They described a simple sentence length approach to sentence alignment and a hybrid multi-feature approach to perform word alignment. They use regression techniques in order to learn parameters which characterize the relationship between the lengths of two sentences in parallel text. They used a multi-feature approach with dictionary lookup as a primary technique and other methods such as local word grouping, transliteration similarity (edit-distance) and a nearest aligned neighbors approach to deal with many-to-many word alignments. Their experiments were based on the EMILLE (Enabling Minority Language Engineering) corpus. They obtained 99.09% accuracy for many-to-many sentence alignment and 77% precision and 67.79% recall for many-to-many word alignment as per the results they published.

III. THE ALIGNMENT MODEL

Alignment is a vital issue in the construction and exploitation of parallel corpora. One of the central modeling problems in statistical machine translation (SMT) is alignment between parallel texts. The alignment methodology tries to identify translation equivalence between sentences, words and phrases within sentences. In most literature, alignment methods are either categorized as association or estimation approaches (heuristic and statistical models). Association approaches use string similarity measures, word order heuristics, or co-occurrence measures (e.g. mutual information scores). The major distinction between statistical and heuristic approaches is that statistical approaches are based on well-substantiated probabilistic models while heuristic ones are not. Estimation approaches use probabilities estimated from parallel corpora, inspired from statistical machine translation, where the computation of word alignments is part of the computation of the translation model.

3.1 IBM Alignment Models 1 thru 3

In the systematic review of statistical alignment models, Och and Ney, 2003 [5] describe the fundamental nature of statistical alignment as trying to model the probabilistic relationship between the source language string s, and target language string t, and the alignment a between positions in s and t. The mathematical notations commonly used for statistical alignment models follow.

$$s_j = s_1, \dots, s_j, \dots, s_J$$

$$t_1 = t_1, \dots, t_i, \dots, t_I \quad (1)$$

Malayalam and English sentences s and t , contain a number or tokens, J and I (equation 1). Tokens in sentences s and t can be aligned, correspond to one another. The set of possible alignments is denoted A , and each alignment from j to i (Malayalam to English) is denoted by a_j which holds the index of the corresponding token i in the English sentence (equation 2).

$$A = \{(j,i) : j=1, \dots, J; i=1, \dots, I\}$$

$$j \rightarrow i = a_j$$

$$i = a_j \quad (2)$$

The basic alignment model using the above described notation is given in equation 3.

$$P(t_1^I | s_1^J)$$

$$P(t_1^I, a_1^I | s_1^J)$$

$$P(t_1^I | s_1^J) = \sum_{a_1^I} P(t_1^I, a_1^I | s_1^J) \quad (3)$$

From the basic translation model $P(s_1^J | t_1^I)$, the alignment is included into equation to express the likelihood of a certain alignment mapping one token in sentence f to a token in sentence t , $P(s_1^J, a_1^I | t_1^I)$. If all alignments are considered, the total likelihood should be equal to the basic translation model probability. The model explained above is the IBM Model 1, which is not considering the word positions.

Model 2 :

One problem of Model 1 is that it has no way of differentiating between alignments that align words on the opposite ends of the sentences, from alignments which are closer. Model 2 adds this distinction. Given source and target lengths (l, M), probability that i th target word is connected to j th source word, the distortion probability is given as $D(i | j, l, s)$. The best alignment can be calculated as follows:

$$a_{j=1}^s [i,j,l,M] \arg_i \max D(i | j, l, s) * T(t_i | s_j) \quad (4)$$

Model 3 :

Indian languages such as Malayalam make use of compound words. English Language does not. This difference makes translating between such languages impossible for certain words, the previous models 1 and 2 would not be capable of mapping one Malayalam word into two English words. Model 3 however introduces fertility based alignment, which considers such one to many translations probable. We uniformly assign the reverse distortion probabilities for model-3. Given source and target lengths (l, M), probability that i th target word is connected to j th source word. The best alignment can be calculated as follows:

$F(\emptyset | s) =$ probability that s is aligned with target words.

$$a_{j=1}^s [i,j,l,M] \arg_i \max D(i | j, M, l) * T(t_i | s_j) * D_{rev}(j | i, M, l) * F(\emptyset | s) \quad (5)$$

3.2 Problem Statements and Suggested Solutions

In IBM models based approaches, the problem of word alignment is divided into several different problems.

The first problem : is to find the most likely translations of an SL word, irrespective of positions.

Solution : This part is taken care of by the translation model. This model describes the mathematical relationship between two or more languages. The main thing is to predict whether expressions in different languages have equivalent meanings.

The second problem : is to align positions in the source language (SL) sentence with positions in the target language (TL) sentence.

Solution : This problem is addressed by the distortion model. It takes care of the differences in word orders of the two languages mentioned. A novel metric to measure word order similarity (or difference) between any pair of languages based on word alignments.

The third problem : is to find out how many TL words are generated by one SL word. Note that an SL word may sometimes generate no TL word, or a TL word may be generated by no SL word (NULL insertion).

Solution : The fertility model is supposed to account for this.

IV. PROPOSED ALIGNMENT MODEL

The proposed system is a combination of corpus based ap-proach and dictionary lookup approach. Alignment step uses corpus based approach first and then dictionary lookup approach. If the corpus has not enough data, the system uses dictionary lookup approach. The approaches are explained below.

4.1 Corpus based Approach

The corpus based approach is based on the first three IBM models and Expectation Maximization (EM) algorithm. The Expectation Maximization (EM) algorithm is used to iteratively estimate alignment model probabilities according to the likelihood of the model on a parallel corpus. In the Expectation step, alignment probabilities are computed from the model parameters and in the Maximization step, parameter values are re-estimated based on the alignment probabilities and the corpus. The iterative process is started by initializing parameter values with uniform probabilities for IBM Model 1. The EM algorithm is

only guaranteed to find a local maximum which makes the result depend on the starting point of the estimation process. This system is implemented using EM algorithm and deals with problem statements. The iterative EM algorithm corresponding to the translation problem can be described as:

Step-1 : Collect all word types from the source and target corpora. For each source word s collect all target words t that co-occurs at least once with s .

Step-2 : Initialize the translation parameter uniformly (uniform probability distribution), i.e., any target word probably can be the translation of a source word t . In this step, there are two main tasks for aligning the source and target sentences. The detailed algorithm of each task is shown in Figure 1 and Figure 2. The first task is pre-processing and the second task is the usage of the first three IBM models.

Pre-processing Phase :

Accept Source Sentence;

Accept Target Sentence;

Remove Stop Word in Source Words (S)

For each Source Sentence S do

Separate into words;

Store Source Words Indexes;

End For

For each Target Sentence T do

Separate into words;

Store Target Words Indexes;

End For

Fig. 1 : Algorithm for Pre-processing

Step - 1: Collect all word types from the source and target corpora.

For each source word s collect all target words t that co occurs at least once with s .

Step - 2 : Any target word (t) probably can be the translation of a source word (s) and the lengths of the source and target sentences are S and T , respectively.

Initialize the expected translation count T_c and Total to 0

Step - 3 : Iteratively refine the translation probabilities.

For $i=1$ to s do

Source Words with N -grams Method

Select Target Words FROM Bilingual corpus (English-Malayalam) WHERE Source Similar si

$total += T(si)$ in corpus

For $j=1$ to t do

If t_j Found in Corpus

$T_c(t_j/si) += T(t_j/si)$

Store Source Word Index and Target Word Index

Align Source Word and Target Word and Store in Corpus

Else if

Use the English Pattern (combine English words with N -grams method)

If $T(si)$ with Target Word found in Corpus

$T_c(t_j/si) += T(t_j/si)$

Store Source Word Index and Target Word Index

Align Source Word and Target Word and Store in Corpus

Else English Word with Null insertion

End If

End For

Calculate Probability T

End For

Fig. 2 : The first three IBM based models algorithm

4.2 Dictionary Lookup Approach

We have used a bilingual Malayalam - English dictionary which consists of 9000 word to word translations. The dictionary lookup approach for alignment is as follows:

Step 1: Find a Root of Malayalam word and English POS Search in Dictionary

Step 2: If not found, we make as unaligned. If yes, Return English Root Word(s) from Dictionary

Step 3: If a match is found, Align Malayalam word with English, O.w, bigram English word.

V. OVERVIEW OF SYSTEM

This system consists of the following steps:

Step 1: Accept a pair of Malayalam and English sentences

Step 2: English is well-developed, and there are many freely available resources for that language. English sentence is passed to a Parser and it will produce Part-of-speech tagged output and root word output.

Step 3: Segment the words in Malayalam sentence using

algorithm in [17], and remove the stop words. In this step, Malayalam sentence is morphological rich. After that, using Tri-Grams method, do the morphological analysis. Each sentence is calculated backward.

Step 4: The output from Step 2 and Step 3 are aligned based on the first three IBM models and EM algorithm using parallel corpus. The result from this step is the aligned words. The high probability words are taken to insert to Parallel Corpus.

Step 5. After Step 4, the remaining unaligned words are aligned using Malayalam-English bilingual dictionary. The lookup approach uses Malayalam root word and English POS in the dictionary to get the English word. Parallel corpus is used as training data set and also the output of the system.

VI. EXPERIMENTAL RESULT

This system uses the Malayalam-English corpus (950 sentence pairs) and 225 sentence pairs for testing. The sentences were at least 4 words long. Malayalam is a very difficult language so that obtaining accurate results is highly impossible. Building a corpus is also very difficult. We declare the expected performance of our alignment Models in terms of precision and recall defined as:

$$\text{Recall} = \frac{W_{\text{correct}}}{W_{\text{Dtotal}}} \times 100\%$$

$$\text{Precision} = \frac{W_{\text{correct}}}{W_{\text{Gtotal}}} \times 100\%$$

$$F\text{-measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \times 100\%$$

Where, Wcorrect is the number of correctly aligned words, WDTotal is the number of words and WSTotal is the number of aligned words by the system. The experimental results are shown in Table 1. By using combination of Corpus based approach and dictionary lookup approach, the precision increased.

Experiment :

S1 is Corpus Based Approach

S2 is Dictionary Lookup Approach

S3 is Corpus Based Approach + Dictionary Lookup Approach

Table 1: Results for experiment

Experiment	S1	S2	S3
Precision(%)	85	87	91

Recall(%)	76	78	82
F-measure(%)	79	82	86

VII. CONCLUSION AND FUTURE WORK

We have shown that building Malayalam-English parallel corpus can be improved by a combination of corpus based approach and dictionary lookup approach. Malayalam languages are morphologically rich. Thus, in future, the proposed model can be enhanced to give better result by using a list of cognates and morphological analysis. This system can be extended as phrase alignment model. We have to work on many to many word alignments and have to test the algorithm for large bilingual corpora.

REFERENCES

- [1] Bing Xiang, Yonggang Deng, and Bowen Zhou, —Diversify and Combine: Improving Word Alignment for Machine Translation on Low-Resource Languages , Proceedings of the ACL 2010 Conference Short Papers, 2010, pages 22–26.
- [2] C. Callison-Burch, D. Talbot, and M. Osborne, —Statistical Machine Translation with Word- and Sentence-Aligned Parallel Corpora . In Proceedings of ACL, Barcelona, Spain, July 2004, pages 175–182.
- [3] D. Wu. —Aligning a Parallel English-Chinese Corpus Statistically with Lexical Criteria In: Proc. of the 32nd Annual Conference of the ACL: 80-87. Las Cruces, NM in 1994. <http://acl.ldc.upenn.edu/P/P94/P94-1012.pdf>
- [4] Eknath Venkataramani and Deepa Gupta, —English-Hindi Automatic Word Alignment with Scarce Resources”. In International Conference on Asian Language Processing, IEEE, 2010.
- [5] F. Och and H. Ney, —A Systematic Comparison of Various Statistical Alignment Models . Computational Linguistics, 29(1):19–52, 2003.
- [6] G. Chinnappa and Anil Kumar Singh, —A Java Implementation of an Extended Word Alignment Algorithm Based on the IBM Models , In Proceedings of the 3rd Indian International Conference on Artificial Intelligence, Pune, India. 2007.
- [7] Helen Langone, Benjamin R. Haskell, Geroge, A. Miller, —Annotating WordNet , In Proceedings of the Workshop Frontiers in Corpus Annotation at HLT-NAACL, 2004.

- [8] Ittycheriah and S. Roukos, —A Maximum Entropy Word Aligner for Arabic-English Machine Translation”. In Proceedings of HLT-EMNLP. Vancouver, Canada, 2005, Pages 89–96.
- [9] J.Martin, R. Mihalcea, and T. Pedersen, —Word Alignment for Languages with Scarce Resources . In Proceedings of the ACL Workshop on Building and Using Parallel Texts. Ann Arbor, USA ,2005, Pages 65–74,.
- [10] Jamie Brunning, Adria de Gispert and William Byrne, —Context-Dependent Alignment Models for Statistical Machine Translation . The 2009 Annual Conference of the North American Chapter of the ACL, pages110–118,Boulder, Colorado, June 2009.
- [11] Li and Chengqing Zong, —Word Reordering Alignment for Combination of Statistical Machine Translation Systems , IEEE, 2008.
- [12] Niraj Aswani and Rpbert Gaizauskas, —A hybrid approach to align sentences and words in English-Hindi parallel corpora . In Proceedings of the ACL Workshop on Building and Using Parallel Texts, June, 2005, page 57-64.
- [13] Pascale Fung and Kenneth Ward Church, — Kvec: A New Approach for Aligning Parallel Texts . In Proceedings of the 15th conference on Computational linguistics. Kyoto, Japan, 1994, Pages 1096-1102.
- [14] P. Koehn, F. J. Och, and D. Marcu, —Statistical Phrase based Translation . In Proceedings of HLT-NAACL. Edmonton, Canada. 2003 ,Pages 81–88.
- [15] P. F. Brown, S. A. Della Pietra, V. J. Della Pietra, and R. L.Mercer, —The Mathematics of Statistical Machine Translation: Parameter Estimation . Computational Linguistics, 19(2):263–311, 1993.
- [16] R. Mihalcea and T. Pedersen, —An evaluation exercise for word alignment . In Proceedings of HLT-NAACL Workshop on Building and Using Parallel Texts: Data Driven Machine Translation and Beyond. Edmonton, Canada., 2003, Pages 1–6.
- [17] W.P.Pa,N.L.Thein, "Disambiguation in Myanmar Word Segmentation",ICCA,February,2009.
- [18] Khin Thandar Nwet, " Building Bilingual Corpus based on Hybrid Approach for Myanmar-English Machine Translation", IJSER, September, 2011.

